

# A Reverse Write Assist Circuit for SRAM Dynamic Write $V_{\text{MIN}}$ Tracking using Canary SRAMs

Arijit Banerjee<sup>1</sup>, Mahmut E. Sinangil<sup>2</sup>, John Poulton<sup>3</sup>, C. Thomas Gray<sup>3</sup>, Benton H. Calhoun<sup>1</sup>

<sup>1</sup>Dept. of ECE, University of Virginia, Charlottesville, VA 22904, USA

<sup>2</sup>NVIDIA, 2 Technology Park Drive, Floor 3, Westford, MA 01886, USA

<sup>3</sup>NVIDIA, 2700 Meridian Pkwy, Suite 100, Durham, NC 27713, USA

<sup>1</sup>E-mail: {ab9ca, bhc2b}@virginia.edu <sup>2,3</sup>E-mail: {msinangil, tgray, jpoulton}@nvidia.com

## Abstract

SRAMs occupy a large amount of area in modern system on chip circuits. With the growing trend of device scaling in deep sub-micron technologies, the 6T SRAM write operation is more vulnerable than the read operation from a failure standpoint. In order to make the SRAMs operate correctly, we must design them with some guard band above the minimum operating voltage ( $V_{\text{MIN}}$ ) by designing for the worst case. In this paper, we investigate a reverse write assist circuit scheme that enables the tracking of SRAM write  $V_{\text{MIN}}$  by using canary SRAM bitcells to track dynamic voltage, temperature fluctuations and aging effects. This circuit ultimately allows us to lower the write  $V_{\text{MIN}}$  below the worst case corner (SF\_85C)  $V_{\text{MIN}}$ , which saves a minimum of 30.7% energy per cycle at the SS\_85C, and a maximum of 51.5% energy per cycle at the FS\_85C corner.

## Keywords

Canary, SRAM, reverse write assist, dynamic  $V_{\text{MIN}}$ .

## 1. Introduction

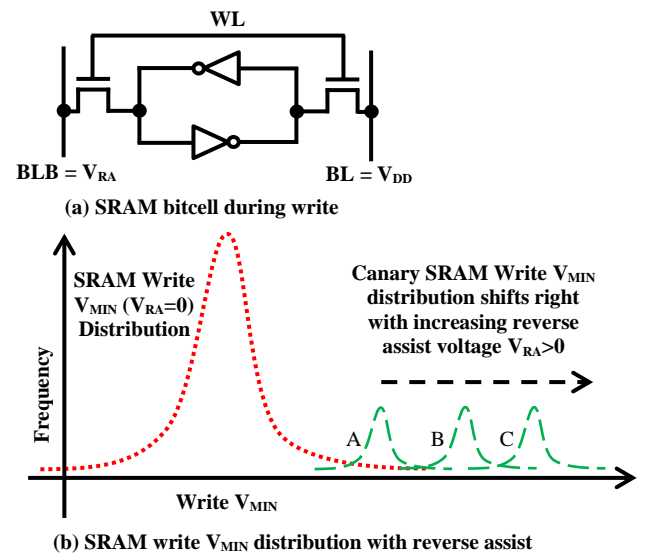
SRAM energy varies quadratically with the supply voltage, so lowering supply voltage lowers energy. There are various ways to lower SRAM supply voltage to lower energy, such as dynamic voltage and frequency scaling (DVFS), using dual rail design for SRAM etc. DVFS is widely used in system on chips (SOCs) to lower the energy [1][2][3] by adjusting the supply voltage and frequency from time to time as required, and SOC level design cost for DVFS is excluded from the SRAM design cost. On the other hand, dual rail [4] designs can be used for energy savings and avoiding readability issues by keeping SRAMs on a higher supply; however, this technique is complicated to implement, and increases design cost in SRAMs and area cost for SOCs. In spite of aforementioned voltage lowering techniques, the SRAM minimum operation voltage ( $V_{\text{MIN}}$ ) poses a bottleneck for SRAM voltage scaling. The SRAM  $V_{\text{MIN}}$  is a function of operating frequency, and it is hard to predict in a real design. So, we design with voltage and timing guard bands. On the other hand, local and global variation make scaling down SRAM  $V_{\text{MIN}}$  more challenging than for logic [5][6], and existing research work shows that SRAM write failure will increase during further scaling [6]. One solution for SRAM read and write  $V_{\text{MIN}}$  improvement is to use assist circuits, such as wordline boosting [7][8][9], negative bitline [7][8][9][10],  $V_{\text{DD}}$  lowering [8][9],  $V_{\text{SS}}$  raising [8][9] etc. for write improvement, and wordline under drive [9], partially suppressed wordline [10],  $V_{\text{DD}}$  boosting [9], negative  $V_{\text{SS}}$  [9] etc. for read improvement.

Assist methods require extra silicon area and power consumption, but can allow for significantly lower SRAM  $V_{\text{MIN}}$ . With time, SRAM circuits age [11][12][13] like all other circuits, and the  $V_{\text{MIN}}$  gets higher and higher, which further adds to the margin necessary for a worst case design [14][15][16].

Hence, predicting the  $V_{\text{MIN}}$  by detecting failures during DVFS can allow corrections to address functionality problems. So, a closed loop solution is ideally required to turn off or on assists or to dynamically adjust the assist voltage when required. Closed loop control can also track the effect of voltage and temperature fluctuations. Hence, there is a need to detect read or write failure dynamically. In this paper, we investigate the use of canary cells to detect failure and to track  $V_{\text{MIN}}$ .

The idea of canary circuits has been studied widely in different fields in circuits [17][18][19]. In SRAMs, the use of canary circuits has been studied by Wang and Calhoun [19][20][21] for predicting the data retention voltage (DRV) during standby, but canaries has not been presented in depth for write or read  $V_{\text{MIN}}$  detection. This paper mainly focuses on the study of canary SRAMs for dynamic write  $V_{\text{MIN}}$ , as device scaling makes a write failure more probable than a read failure [6].

In this paper, Section 2 discusses assists and reverse assists. In Section 3, we discuss the effect of reverse assist in



**Figure 1: SRAM write operation using bitline type reverse assist and write  $V_{\text{MIN}}$  distributions with reverse assist (A, B, C's are canary  $V_{\text{MIN}}$  distributions)**

canary SRAMs. We develop a methodology based on probability theory and use the concept of canary SRAM to quantify the output metrics in Section 4. Section 5 gives simulation results using our methodology. Section 6 describes the circuit implementation, and we propose a canary SRAM architecture using BL type reverse assist and an algorithm to track SRAM  $V_{\text{MIN}}$  with canary reverse assist in Section 7. We discuss the power and area tradeoffs in Section 8 and conclude in Section 9.

## 2. Peripheral assist methods and reverse assists

There are many ways to create canary circuits in SRAMs. One method is to modify the SRAM core bitcell to fail earlier than a population of SRAM bitcell during the read or write operation. We can have built in control in the canary bitcell to tune the canary to change the failure point. However, this type of canary bitcell may not track same as core bitcells over variation. Another option is to use a shorter wordline pulse width modulator circuit for canaries to make the write/read operation more difficult to fail earlier. In order to get a precisely controlled wordline pulse width, extra wordline delay control circuit is required, which will increase the area overhead in SRAM decoder and may cause abutting problems in layout.

An assist in the SRAM context means an auxiliary circuit that helps improve write-ability [7][8][9][10][4], readability [9][10], or read stability [4]. We define a reverse assist as an auxiliary circuit that degrades the write-ability or readability of an SRAM cell. In this work, we use the same core SRAM bitcell as a canary SRAM, but we apply a reverse assist to degrade the canary SRAM bitcell write-ability or readability.

The advantage of a reverse assist for a canary SRAM is to use the same SRAM core bitcells as canaries to track the core cells better. Also, a user can control the reverse assist with low overhead to fine tune the failure point of the canary SRAM bitcells dynamically.

In the context of this paper, assist or reverse assist will always refer to a bitline (BL) type assist or reverse assist. In core SRAMs without any assist, either BL or BLB is pulled down (Figure 1 (a)) to  $V_{\text{RA}}=0\text{V}$  during a write, while the other node (BLB or BL) is kept at  $V_{\text{DD}}$ . Usually a BL type assist [7][8][9][10] is used to improve the dynamic write-ability of the SRAM bitcells by pulling the bitline or bitline bar (BLB) node below the ground voltage ( $V_{\text{SS}}$ ). On the other hand, in canary SRAMs a reverse assist will pull the BL/BLB node to a positive voltage, say for example  $V_{\text{RA}}=0.1\text{V}$ , while the other BLB/BL is kept at  $V_{\text{DD}}$  as shown in Figure 1 (a). This will degrade the dynamic write-ability of the canary SRAM bitcells.

## 3. Effect of reverse assist on canary SRAMs and canary design metrics

The ability to write in an SRAM bitcell is called bitcell write-ability. There are two widely used metric for write-ability as write static noise margin (WSNM) known as the static write-ability metric, and another metric called critical wordline pulse width for write ( $T_{\text{CRIT}}$ ) known as the dynamic write-ability metric for SRAMs. WSNM assumes the wordline pulse width to be infinite which overestimates the

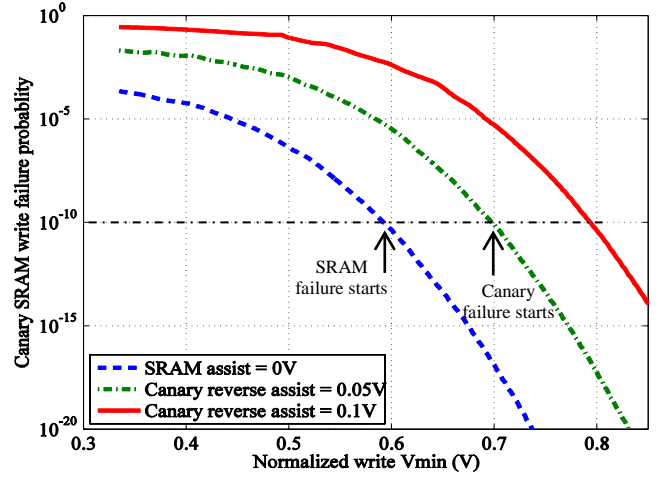


Figure 2: Canary SRAM write failure probability vs. normalized write  $V_{\text{MIN}}$

Table 1: Input and output design metrics for canary SRAM design

Input Metrics	
N	Number of SRAM bits on a chip
$Y_{\text{SRAM}}$	Core SRAM target yield
C	Number of canary SRAM bits
$F_{\text{th}}$	Canary failure threshold condition
$V_{\text{RA}}$	Canary BL type reverse assist voltage
Output Metrics	
$P_{\text{fc}}$	Canary SRAM chip failure probability

static write-ability metric, but  $T_{\text{CRIT}}$  assumes a realistic finite wordline pulse width as SRAM write operation is a dynamic process. Using a write assist in SRAMs causes the spread of the distribution of  $T_{\text{CRIT}}$  to decrease and to shift the  $V_{\text{MIN}}$  to a lower value [8]. Hence, applying a reverse assist to the canary bitcells, relative to the actual SRAM cells, the canary write  $V_{\text{MIN}}$  distribution ‘A’ will shift to a higher  $V_{\text{MIN}}$  distribution ‘B’ or ‘C’ as shown in the Figure 1 (b). Hence, by applying a reverse assist to the canary SRAM bitcells, the  $V_{\text{MIN}}$  of the canary SRAM bitcells will increase to cause canary failures earlier than the core SRAM bitcells.

Our goal is for the canary SRAMs to start to fail before a single failure in a given number of SRAM bits, say a million or billion bits. Figure 2 shows the simulated dynamic write probability of failure ( $P_{\text{fail}}$ ) vs. write  $V_{\text{MIN}}$  plots for core SRAM cells and canary SRAM cells with varying degrees of reverse assist. Here, we use the extracted 6T bitcell netlist with the setup shown in Figure 1 (a) and simulate transient write operation using a commercial 28nm technology with HSPICE using an importance sampling algorithm [5][22][23][24] to get the  $P_{\text{fail}}$  vs. dynamic write  $V_{\text{MIN}}$  data. We can see that for the same probability of failure, the canary SRAM using reverse assist has a higher write  $V_{\text{MIN}}$  than the core SRAMs without any assist.

Table 1 defines the input and output design metrics for canary SRAM design. A write failure probability for core SRAM corresponds to a certain number of SRAM bits  $N$  on a chip with a target yield  $Y_{SRAM}$ . On the other hand, tracking dynamic write failure of actual SRAM bits requires a certain number of canary bits  $C$ . Other important input knobs are the canary failure threshold condition ( $F_{th}$ ) and reverse assist voltage ( $V_{RA}$ ).  $F_{th}$  condition is the number of canary cells allowed to fail before one in  $N$  SRAM core bits fails. For example, if a user defines  $F_{th}=8$  for  $C=32$  canaries in a chip, then an action can be taken if 8 canaries fail to write out of 32 canaries. The possible actions are either turning on assists for the actual SRAM bits or stopping voltage and frequency scaling further. Also, the amount of degradation of writeability in canaries can be controlled by tuning the  $V_{RA}$ . The two input metric knobs available post-fabrication to the user are  $V_{RA}$  and  $F_{th}$  for canary SRAMs. All other input metrics are set at design time. For the output metric, we define  $P_{fc}$  as *canary SRAM chip failure probability*, which is the probability that the canary SRAM bitcells will be unable to fail earlier than one in  $N$  SRAM core bitcells. For example, if  $P_{fc}=10^{-6}$  for a given  $N=10^7$  SRAM bits with  $Y_{SRAM}=99\%$ ,  $C=32$ , and  $F_{th}=1$ , then the canary chip failure probability will denote that in one in a million 10Mb chips, the 32 canary cells will *not* experience a single bit failure prior to the first failure from ten million SRAM bits on the chip.

Now, for core SRAMs, if the bit failure probability in a write operation is given by  $P_f$ , then the core SRAM bitcell success probability is given by  $P = (1 - P_f)$ , the SRAM chip success probability is given by  $P_{chip} = P^N$ , and the SRAM chip failure probability is given by  $P_{f_{chip}} = (1 - P_{chip})$ . Now, the SRAM chip yield for 'k' or less chip failures in total of 'J' chips can be given by:

$$Y_{SRAM(J,k)} = \sum_{i=0}^{i=k} P_{chip}^{(J-i)} * (1 - P_{chip})^i * \binom{J}{i} \quad (1)$$

From (1), for a given value of  $Y_{SRAM}$  and  $N$ , we can calculate the corresponding SRAM bit failure probability  $P_f$  for write failures. Similarly, if the canary bit failure probability is given by  $p_f$ , then the probability of canary bits being unable to fail earlier than a given number of core SRAM bits with  $C$  number of canary bitcells with  $F_{th}=k$  condition can be given by:

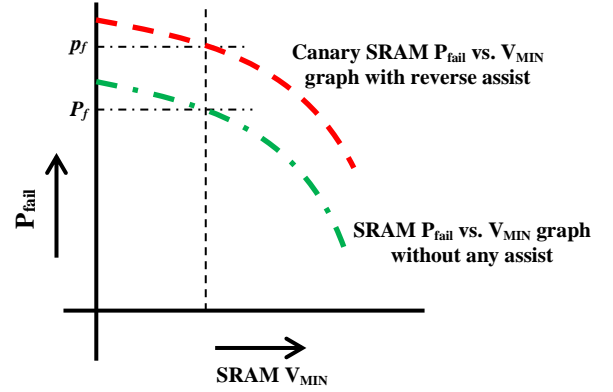
$$P_{fc} = \sum_{i=0}^{i=k} p_f^i * (1 - p_f)^{C-i} * \binom{C}{i} \quad (2)$$

Hence, (1) and (2) relate the input metrics  $N$ ,  $Y_{SRAM}$ ,  $P_f$ ,  $C$ ,  $F_{th}$ , and  $p_f$  to canary chip failure probability  $P_{fc}$ , which is our final output metric for observation.

#### 4. Calculation Methodology for Canary Chip Failure Probability

Figure 3 shows the methodology to calculate the canary bit failure probability  $p_f$  using SRAM bit failure probability  $P_f$ . The plot shows  $P_{fail}$  vs.  $V_{MIN}$  for core SRAM bitcells

without any assist, and canary SRAM  $P_{fail}$  vs.  $V_{MIN}$  with reverse assist. The research question we are addressing here is, for a given value of  $Y_{SRAM}$ , and  $N$ , what is the

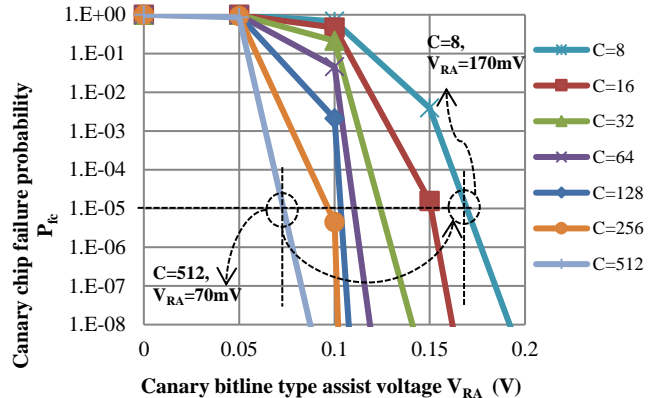


**Figure 3: Methodology to calculate canary chip failure probability**

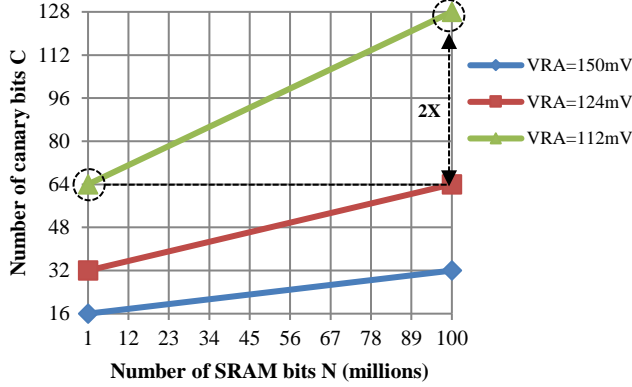
corresponding SRAM bit failure probability  $P_f$ , and what should be the corresponding bit failure probability  $p_f$  for canary SRAM bits? We also want to know how the input metric  $C$  influences the canary chip failure probability  $P_{fc}$ . In order to connect the two equations (1) and (2), we use the setup mentioned in Section 3 to get the  $P_{fail}$  vs.  $V_{MIN}$  data for different reverse assist voltages, which represents the data for the canary SRAMs bitcells. We also got the  $P_{fail}$  vs.  $V_{MIN}$  data for core SRAM 6T bitcells without any assist using the same simulation setup. First of all, we calculate the corresponding  $P_{fail}$   $P_f$  for SRAM bitcells using (1), and then we calculate the corresponding  $V_{MIN}$  for the SRAM bitcells using the  $P_{fail}$  vs.  $V_{MIN}$  simulated data (Figure 3). Then we calculate the corresponding  $P_{fail}$   $p_f$  for canary SRAM bitcells with same  $V_{MIN}$  obtained from the canary SRAM  $P_{fail}$  vs.  $V_{MIN}$  simulated data, which is shown in Figure 3. Finally, we plug the value of  $p_f$  into (2), and calculate the corresponding canary chip failure probability  $P_{fc}$ .

#### 5. Simulation Results for Canary SRAM Chip Failure Probability

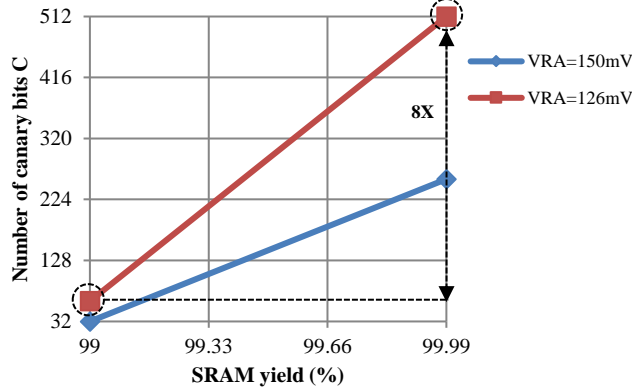
In order to get the trends of the input metric vs. the output metric variation, we use the calculation method described in Section 4 and calculate the output metric for reverse assist voltages of  $V_{RA}=0V, 0.05V, 0.1V, 0.15V$  and



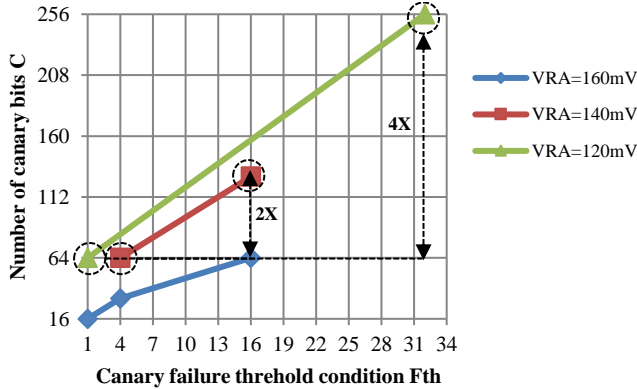
**Figure 4: Canary chip failure probability vs. reverse assist voltage for 1 million SRAM bitcells with 95% yield @ TT\_85C**



**Figure 5: Trend for C vs. N with 95% SRAM yield at constant  $P_{fc}=10^{-5}$  for different  $V_{RA}$  voltages @ TT\_85**



**Figure 6: Trend of C vs.  $Y_{SRAM}$  with 100 million SRAM bitcell at constant  $P_{fc}=10^{-5}$  for different  $V_{RA}$  voltages @ TT\_85C**



**Figure 7: Trend of C vs.  $F_{th}$  with 100 million SRAM bitcell at constant  $P_{fc}=10^{-5}$  for different  $V_{RA}$  voltages @ TT\_85C**

0.2V. Figure 4 shows that same canary chip failure probability  $P_{fc}=10^{-5}$  can be achieved by either increasing the number of canaries to  $C=512$  with a lower  $V_{RA}=70mV$  or decreasing  $C=8$  with a higher  $V_{RA}=170mV$ . In order to get a trend of C vs. N, C vs.  $Y_{SRAM}$ , and C vs.  $F_{th}$  for a constant

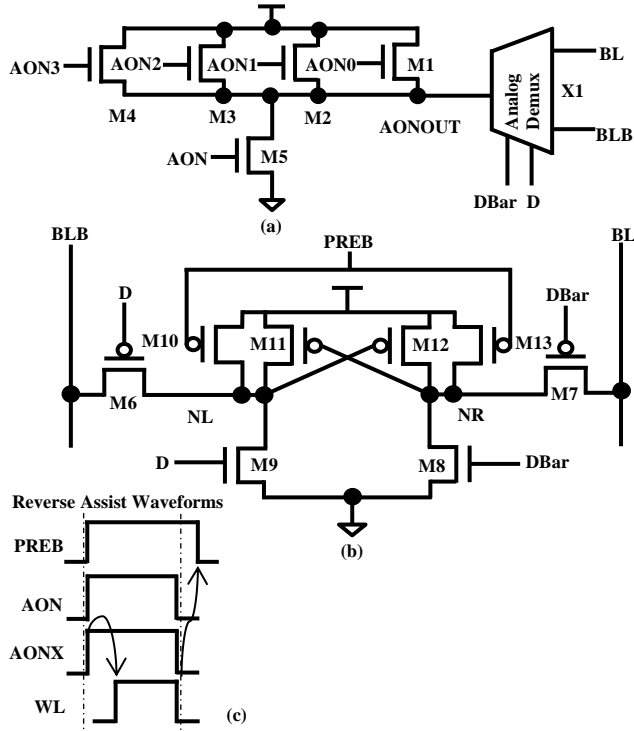
$P_{fc}=10^{-5}$  with different values of  $V_{RA}$  we interpolated the data for  $V_{RA}$  in between known  $V_{RA}$  values.

Figure 5 shows the trend of the number of canary bits C vs. the number of SRAM bits N. We can see that increasing N two orders of magnitude from 1 million to 100 million bits, causes the number of canaries C required to maintain the same canary chip failure probability of  $P_{fc}=10^{-5}$  (at the different reverse assist voltages) to double. Figure 6 show trend of number of canary bits C vs. SRAM yield  $Y_{SRAM}$ . We can see that in order to keep the same canary chip failure probability of  $P_{fc}=10^{-5}$ , while increasing the SRAM yield from 99% to 99.99%, the number of canary bits has to be increased by 8X from  $C=64$  to  $C=512$  for  $V_{RA}=126mV$  BL type reverse assist.

Similarly, Figure 7 shows the trend of the number of canary bits C vs. canary failure threshold condition  $F_{th}$  while keeping other input metrics constant. We can see that to maintain the same canary chip failure probability roughly at  $P_{fc}=10^{-5}$ , increasing the failure threshold from  $F_{th}=4$  to  $F_{th}=16$ , requires 2X more canary cells than that of the  $C=64$  at  $V_{RA}=140mV$  BL type reverse assist voltage. On the other hand, for reverse assist voltage  $V_{RA}=120mV$ , a change of 32X in  $F_{th}$  condition requires 4X increase in C from  $C=64$  to  $C=256$  to maintain the same  $P_{fc}=10^{-5}$ .

## 6. Circuit implementation of BL type reverse assist

We assume that a reverse assist will be integrated inside the existing core SRAM I/O as canary I/O, and therefore it requires additional circuitry. Possible ways of creating a reverse assist is to use a positive charge pump, or an analog closed loop voltage reference, or a voltage divider circuit, etc. to generate the reverse assist voltage for the BL type reverse assist. A charge pump and analog closed loop variable voltage reference would have caused much higher design and area overhead per canary I/O. Also, we found that a PMOS-NMOS voltage divider has much higher variation of the output voltage than an NMOS-NMOS



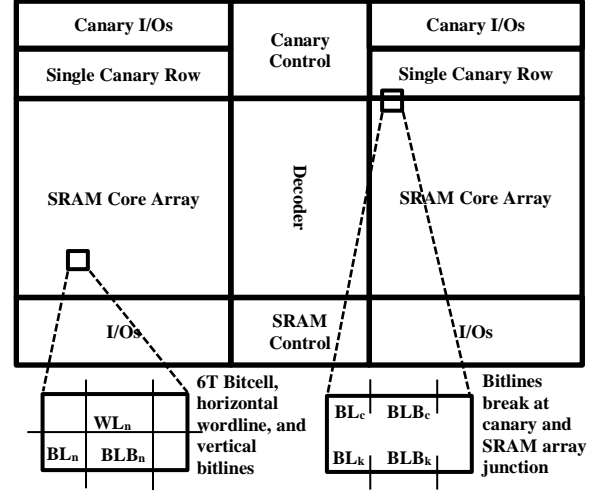
**Figure 8: Canary SRAM reverse assist circuit with write driver and reverse assist waveforms**

voltage divider. To write canaries, we propose a novel reverse assist (Figure 8(a)) which generates the positive bias voltage ( $V_{RA}$ ) for BL/BLB, and a write driver (Figure 8(b)) which pulls up the other node BLB/BL to  $V_{DD}$ . Here, signals AON0, AON1, AON2, and AON3 are cumulatively represented by the name AONX in Figure 8(c). Signal AON creates the  $V_{RA}$  using an NMOS-NMOS voltage divider circuit by selecting M5 and M1-M4. Here, the  $V_{RA}$  of the voltage divider at node AONOUT is channeled to BL/BLB using an analog de-multiplexer X1. During a write operation using reverse assist, D or DBar turns on M9/M8 to pull down one of the NL/NR nodes to ground (Figure 8(b)). This pulls up NR/NL accordingly through cross coupled M12 and M11. However, only the pulled up node NR/NL gets connected to desired BLB/BL nodes by M7 or M6. Thus, M6 and M7 separate the internal pulled down node NL/NR from BL/BLB by turning off M6 or M7, which allows us to connect the reverse assist voltage node AONOUT to BL/BLB node using the analog de-multiplexer X1. For this experiment, we propose to size the analog demultiplexer, M1-M5 sufficiently to discharge the BL/BLB and to generate a minimum of 50mV and maximum of 200mV of reverse assist in a write operation.

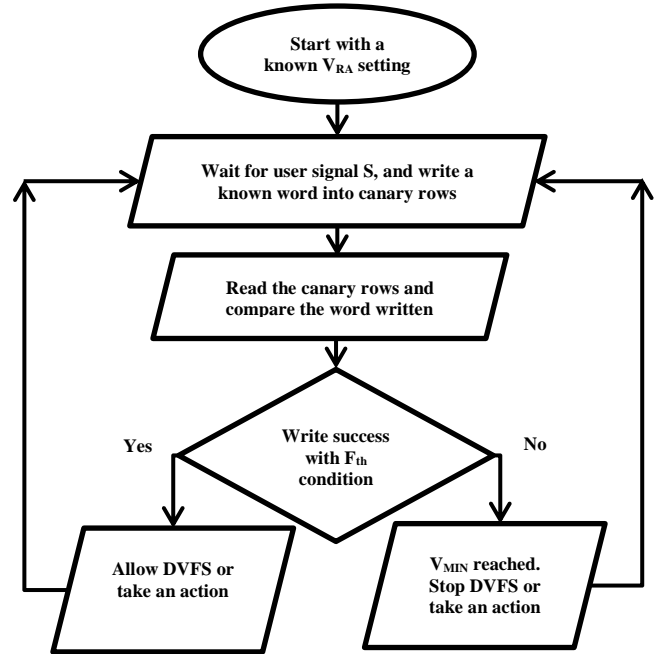
## 7. Block diagram of canary SRAM architecture and an algorithm to track SRAM $V_{MIN}$

In order to implement the circuit proposed in Section 6, we propose a canary SRAM architecture and an algorithm to track SRAM  $V_{MIN}$  in this Section. The block diagram of the proposed canary architecture is shown in Figure 9. In this block diagram, the canary I/Os, canary control, and single canary bitcell rows constitute the canary SRAM. This canary block can adjoin a core SRAM macro as shown in

Figure 9, which has two SRAM core arrays, a decoder, I/Os, and SRAM control logic. In this case, the canary control can directly talk to the SRAM control logic. The wordlines are oriented horizontally, and bitlines vertically in the SRAM core array and in canary bitcell rows. In order to operate the canaries independently of SRAM, the bitlines break at the junction of SRAM core array and canary row junction. Also, the canary can sit distantly from the SRAM macros. In



**Figure 9: Block diagram of canary SRAM inside SRAM macro (not in scale)**



**Figure 10: Canary reverse assist control and  $V_{MIN}$  tracking algorithm**

the first case, if the canary SRAM is integrated in all the SRAM macros, it can track local and global voltage, frequency, temperature fluctuations on the power grid, variation in corners, and aging effects in a large SOC. On the other hand, a standalone single canary SRAM macro can only track the global variation effects in corners, aging etc. in an SOC. The reverse assist circuit described in Figure 8 sits inside each individual canary I/O, and to reduce the



effect of local variation, the AONOUT (Figure 8(a)) signal is shared among the canary I/Os.

Figure 10 shows our proposed algorithm for tracking the SRAM  $V_{\text{MIN}}$ . Initially, the canary control state machine (CCSM) starts with a known setting of  $V_{\text{RA}}$  during boot up. This known  $V_{\text{RA}}$  setting corresponds to the SRAM  $V_{\text{MIN}}$  at a certain process corner with a specific SRAM size of  $N$  bits,  $C$  number of integrated canaries in the SRAM for a constant  $P_{fc}$  (Figure 4) etc. parameters. After applying the initial  $V_{\text{RA}}$  setting, the canary state machine waits for a user signal 'S.' If the user allows canary operation by setting signal 'S,' then the CCSM writes a known word into the canary rows in the first cycle and reads it back from the canary rows to compare with the existing known word value in the second cycle. Word matching with less than or equal to  $F_{\text{th}}$  number of canary failures signifies a successful write in canaries in the previous cycle, else write fails. Write failure in canaries with an  $F_{\text{th}}$  condition indicates an imminent SRAM failure. In this situation, the CCSM can signal the DVFS control logic in SOC to stop voltage scaling further, or take a user defined action like stalling the memory access for couple of cycles or turn on assist in SRAMs, etc. Otherwise, further voltage scaling is allowed or a user defined action like turning off the SRAM assist etc. can be taken. Thus, this algorithm can track the  $V_{\text{MIN}}$  of each individual SRAM macros with built in canaries. Also, the CCSM can quantify the number of failures and use this value to set the  $F_{\text{th}}$  value. Moreover, a user can update the initial  $V_{\text{RA}}$  setting using on-chip temperature or aging sensor data and simulation data of  $P_{fc}$  to track the write  $V_{\text{MIN}}$  more precisely.

## 8. Power and area tradeoff for the BL type reverse assist circuit with write driver

All the power (total of dynamic and leakage) and area tradeoff numbers related to  $P_{fc}$  are calculated with the assumption that the total number of SRAM bits  $N$  is 100 million in an SOC with SRAM yield  $Y_{\text{SRAM}}$  of 99%, elsewhere tradeoff numbers are calculated using some assumption of wordline driver width, I/O height, average bitline energy and bitcell energy per bit etc. parameters. A single canary SRAM, whether or not integrated inside a core SRAM macro, will not be able to track local fluctuations of voltage and frequency in power bus, and temperature in all 100 million core SRAM bits. This is because those bits are distributed all over the SOC, and the voltage etc. fluctuations will vary from place to place in each macros. As this total  $N$  number of SRAM bits can be divided into  $M$  number equal or unequal sized SRAM macros, we need to quantify the effect of area and energy overhead of canaries vs. average size of SRAM macros.

If the canary SRAM is integrated inside a core SRAM macro, the number of canary I/Os, assuming column mux

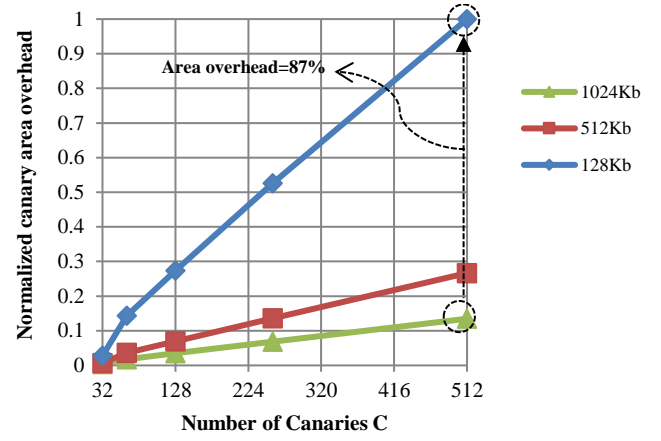


Figure 11: Canary area overhead normalized vs. number of canaries for different SRAM sizes

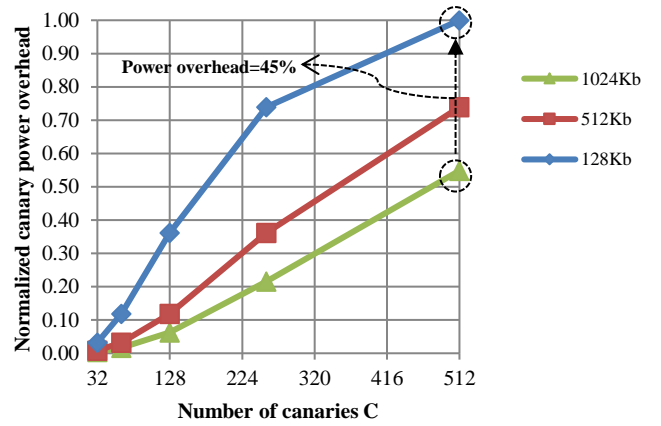


Figure 12: Normalized canary total power overhead vs. number of canaries C with constant  $V_{\text{RA}}=50\text{mV}$  for different SRAM sizes at 1GHz TT\_85C corner

(CM) 4 scenario, required are same as number of SRAM I/Os to make a rectangular shaped symmetric total SRAM macro (Figure 9). Hence,  $C$  is dependent on the number of I/Os in the SRAM. If a designer chooses a logical macro size of 128 words, 64 bits with CM 4 (128x64x4), he has to use  $C=64 \times 4=256$ . Hence, the SRAM size fixes the number of canaries in integrated canary SRAM macros; still we can use standalone canary SRAMs of user defined size in between core SRAM macros. On the other hand, canary I/Os occupy much bigger area compare to the canary bitcells, which dominates the canary SRAM area overhead. Figure 11 shows that increasing number of canary  $C$  increases the area overhead, and with same  $C=512$  canaries (128 I/Os with CM=4) the overhead is 87% more in smaller 128Kb macros than the bigger one of size 1024Kb. Hence, area can be traded off for better tracking of small sized SRAM macros'  $V_{\text{MIN}}$  across an SOC accurately. Figure 12 shows that for the same  $C=512$  number of canaries, the 128Kb SRAM macro has roughly 45% higher total power (dynamic and leakage) overhead than a 1024Kb SRAM macro with the same  $V_{\text{RA}}=50\text{mV}$  at the TT\_85C corner with operating frequency of 1GHz. On the other hand, Figure 13 shows for the change of  $V_{\text{RA}}=50\text{mV}$  to  $V_{\text{RA}}=150\text{mV}$ , the canary power overhead in SRAMs increases by 30% for a 512Kb

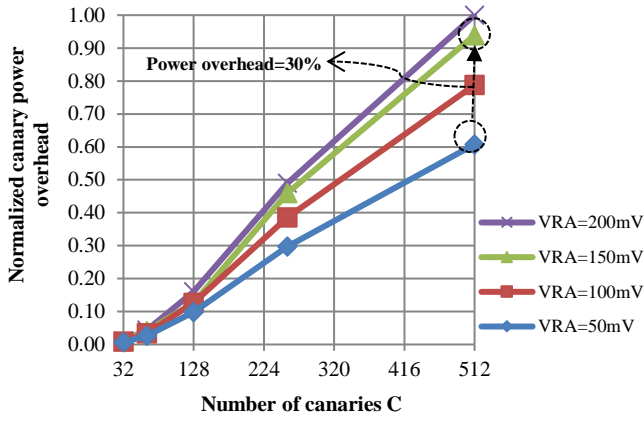


Figure 13: Normalized canary total power overhead vs. number of canaries C with N=512Kb SRAM for different reverse assist voltages at 1GHz TT\_85C corner

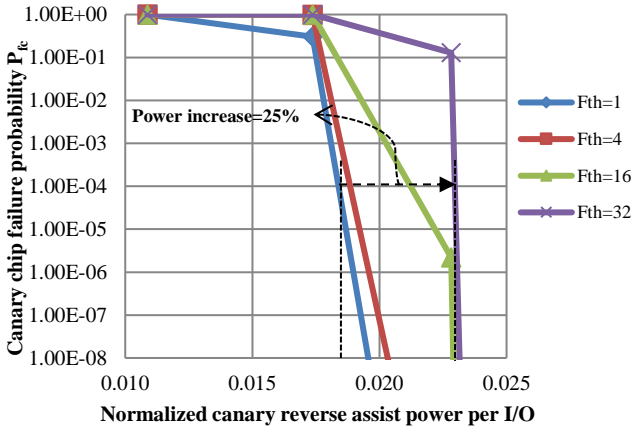


Figure 14: Canary chip failure probability vs. normalized canary reverse assist total power for increasing  $F_{th}$  conditions at 1GHz TT\_85C corner

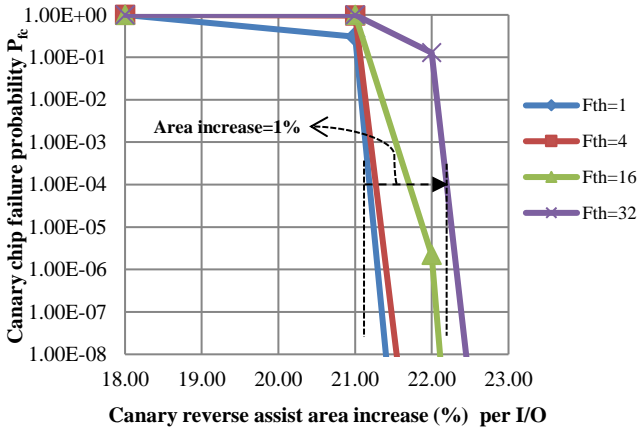


Figure 15: Canary chip failure probability vs. canary reverse assist area increase per I/O for increasing canary failure threshold condition

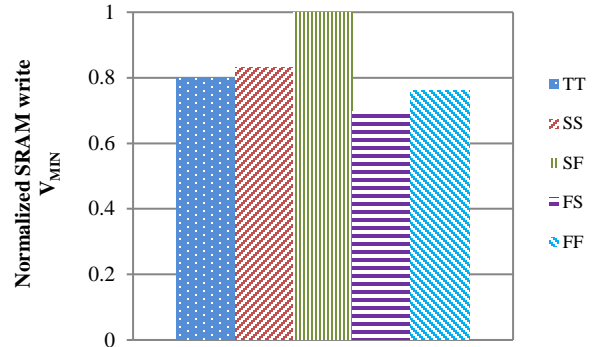


Figure 16: Normalized SRAM write  $V_{MIN}$  for 100 million SRAM bits with 99% yield constraints at TT\_85C corner

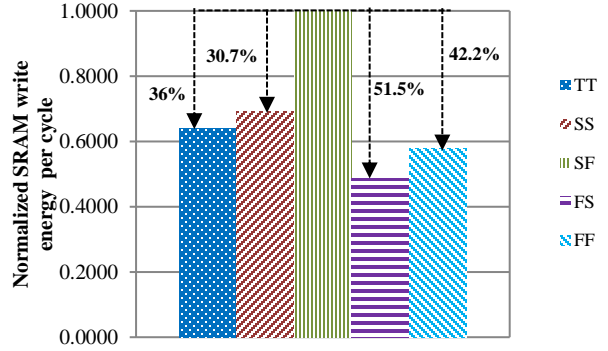


Figure 17: Normalized SRAM write energy per cycle at  $V_{MIN}$  for 100 million SRAM bits with 99% yield constraints at TT\_85C corner

macro with  $C=512$  at 1GHz operating frequency. Figure 14 and Figure 15 show the  $P_{fc}$  vs. power cost and area cost for increasing  $F_{th}$  values. We can see that for the same  $P_{fc}=10^{-4}$  at 1GHz, increasing the  $F_{th}$  condition from  $F_{th}=1$  to  $F_{th}=32$  with  $C=128$  cells using BL type reverse assist, the power cost increases by 25% at TT\_85C, and the canary I/O area cost is increased by 1%.

Ultimately, the normalized SRAM  $V_{MIN}$  for 100 million SRAM bits with 99% SRAM yield at 85C temperature is simulated and shown in Figure 16. As per our simulation results, the canary SRAMs can be used to track the write  $V_{MIN}$  of the SRAM bits with a specified confidence, and the normalized write energy corresponding to the core SRAM  $V_{MIN}$  is shown in Figure 17. We can see that at the TT\_85C corner we can operate SRAMs with 36% lower energy cost than that of the worst case  $V_{MIN}$  in SF\_85C corner, which would set the guard band. The least energy savings can be achieved at the SS\_85C corner as 30.7%. The maximum energy savings from Figure 17 can be found to be in FS\_85C corner, which is 51.5% lower than the worst case. And at FF\_85C corner, the energy savings can reach up to 42.2% with respect to the worst case energy in SF\_85C corner.

## 9. Conclusion

We conclude that the canary SRAM concept using a reverse assist is a promising solution to predict core SRAM failure resulting from write-ability problems. Canary SRAM enables the tracking of SRAM write  $V_{MIN}$  by using reverse

assist to track dynamic voltage, temperature fluctuations and aging effects, and allows us to take necessary actions which can be in the form of turning on assists, stalling the memory access, slowing down operating frequency, or boosting supply voltage. Here, we do all the probability calculations based on importance sampling algorithm. However, choosing an incorrect importance sampling distribution can mispredict the  $V_{\text{MIN}}$  to cause higher energy dissipation or SRAM failures before canaries. The area and power overhead for canary SRAMs are lower for bigger SRAM macros, and they depend on the number of SRAM I/Os in integrated canary SRAM macros. We can qualitatively say that although increasing the canary failure threshold condition  $F_{\text{th}}$  degrades the canary chip failure probability, it also rejects the extreme canary outliers. Moreover, using canaries, SRAM write  $V_{\text{MIN}}$  in different corners can be reduced below the traditional worst case  $V_{\text{MIN}}$ , which saves energy. Finally, we conclude that the canary SRAM for dynamic write  $V_{\text{MIN}}$  tracking works in simulation and theory.

## 10. Acknowledgements

We thank DARPA and NVIDIA for supporting this research work. This work was funded in part by DARPA through a subcontract from NVIDIA. This research was, in part, funded by the U.S. Government. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the U.S. Government.

## 11. References

- [1] S. Herbert and D. Marculescu, "Analysis of dynamic voltage/frequency scaling in chip-multiprocessors", in *Proc. Int. Symp. Low Power Electron. Design*, pp. 38-43, 2007.
- [2] A. Genser et al., "Power emulation based DVFS efficiency investigations for embedded systems", in *Proc. Int. Symp. Syst. Chip (SoC)*, pp. 173-178, 2010.
- [3] R. Airoidi et al., "Improving reconfigurable hardware energy efficiency and robustness via DVFS-scaled homogeneous MPSoC," in *Proc. IEEE International Symposium on Parallel and Distributed Processing Workshops and Phd Forum (IPDPSW)*, pp. 286-289, 2011.
- [4] J. Pille et al., "Implementation of the cell broadband engine in a 65 nm SOI technology featuring dual-supply SRAM arrays supporting 6 GHz at 1.3 V," in *IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, pp. 322-323, 2007.
- [5] B. Zimmer et al., "SRAM assist techniques or operation in a wide voltage range in 28 nm CMOS," in *IEEE Trans. Circuits Syst. II*, vol. 59, no. 12, pp. 853-857, 2012.
- [6] A. Bhavnagarwala et al., "Fluctuation limits & scaling opportunities for CMOS SRAM cells," in *IEDM Tech. Dig.*, 2005, pp. 659-662, 2005.
- [7] N. Shibata et al., "A 0.5-V 25-MHz 1-mW 256-kb MTCMOS/SOI SRAM for solar-power-operated portable personal digital equipment—Sure write operation by using step-down negatively overdriven bitline scheme", in *IEEE J. Solid-State Circuits*, pp. 728-742, 2006.
- [8] V. Chandra et al., "On the efficacy of write-assist techniques in low voltage nanoscale SRAMs," in *Proc. Des. Autom. Test Eur.*, pp. 345-350, 2010.
- [9] R.W. Mann et al., "Limits of Bias Based Assist Methods in Nano-Scale 6T SRAM", in *Proceedings of Quality Electronic Design Symposium*, pp. 1-8, 2010.
- [10] Jonathan Chang et al., "A 20nm 112Mb SRAM in High- $\kappa$  Metal-Gate with Assist Circuitry for Low-Leakage and Low-VMIN Applications" in *Solid-State Circuits Conference Digest of Technical Papers (ISSCC), 2013 IEEE International*, pp. 316 - 317, 2013.
- [11] A. Shah, H. Mahmoodi, "Thermal estimation for accurate estimation of impact of BTI aging effects on nano-scale SRAM circuits," in *SOC Conference (SOCC), 2010 IEEE International*, pp. 230-235, 2010.
- [12] A. Bansal et al., "Impact of NBTI and PBTI in SRAM bit-cells: Relative sensitivities and guidelines for application-specific target stability/performance," in *IEEE IRPS*, pp. 745-749, 2009.
- [13] S. C. Yang et al., "Timing control degradation and NBTI/PBTI tolerant design for write-replica circuit in nanoscale CMOS SRAM", in *Proc. IEEE Int. Symp. VLSI Design, Autom., Test*, pp. 162-165, 2009.
- [14] S. Nalam et al., "Dynamic write limited minimum operating voltage for nanoscale SRAMs," in *Proc. Des. Autom. Test Eur.*, pp. 1-6, 2011.
- [15] J. Wang et al., "Two Fast Methods for Estimating the Minimum Standby Supply Voltage for Large SRAMs", in *Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, vol. 29, issue 12, pp. 1908-1920, 2010.
- [16] E. Karl et al., "A 4.6ghZ 162Mb SRAM design in 22nm trigate CMOS technology with integrated active  $V_{\text{min}}$ -enhancing assist circuitry," in *IEEE Int. Solid-State Circuits Conf. Dig. Tech. Papers*, San Francisco, CA, pp. 230-232, 2012.
- [17] B. H. Calhoun and A. P. Chandrakasan, "Standby power reduction using dynamic voltage scaling and canary flip-flop structures", in *IEEE J. Solid-State Circuits*, vol. 39, no. 9, pp. 1504-1511, 2004.
- [18] Y. Otsuka et al., "Multicore energy reduction utilizing canary FF," in *Communications and Information Technologies (ISCIT), 2010 International Symposium*, pp. 922- 927, 2010.
- [19] J. Wang and B. Calhoun, "Canary replica feedback for near-DRV standby  $V_{\text{DD}}$  scaling in a 90 nm SRAM," in



- Proc. Custom Integrated Circuit Conf. (CICC '07)*, pp. 29–32, 2007.
- [20] J. Wang, and B. H. Calhoun, "Techniques to Extend Canary-based Standby  $V_{DD}$  Scaling for SRAMs to 45nm and Beyond", in *IEEE Journal of Solid-State Circuits*, vol. 43, pp. 2514-2523, 2008.
- [21] J. Wang et al., "An Enhanced Canary-based System with BIST for SRAM Standby Power Reduction", in *Transactions on VLSI Systems (TVLSI)*, pp. 909-914, 2011.
- [22] L. Dolecek et al., "Breaking the simulation barrier: SRAM evaluation through norm minimization," in *Proc. IEEE/ACM Int. Conf. Comput.-Aided Des.*, pp. 322-329, 2008.
- [23] R. Kanj et al., "Mixture importance sampling and its application to the analysis of SRAM designs in the presence of rare failure events," in *Proc. ACM/IEEE Des. Autom. Conf.*, pp. 69-72, 2006.
- [24] A. Singhee et al., "Recursive Statistical Blockade: An Enhanced Technique for Rare Event Simulation with Application to SRAM Circuit Design", in *International Conference on VLSI Design, India*, pp. 131-136, 2008.